

# INTERCHART: Benchmarking Visual Reasoning Across Decomposed and Distributed Chart Information

Anirudh Iyengar Kaniyar Narayana Iyengar <sup>1\*</sup>, Srija Mukhopadhyay <sup>2\*</sup>,  
Adnan Qidwai <sup>2\*</sup>, Shubhankar Singh <sup>3</sup>, Dan Roth <sup>4</sup>, Vivek Gupta <sup>1</sup>

<sup>1</sup>Arizona State University, <sup>2</sup>IIT, Hyderabad, <sup>3</sup>Mercer Mettl, <sup>4</sup>University of Pennsylvania

*\*Equal Contribution*



# When Charts Mislead, Policies Fail ?

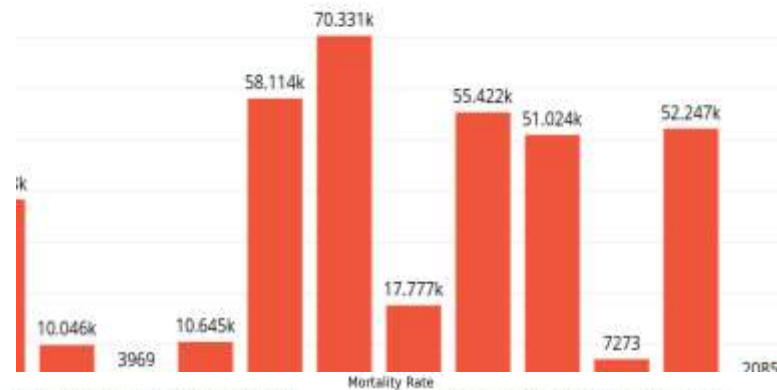
## COVID-19 Data Misrepresented by Georgia Health Department



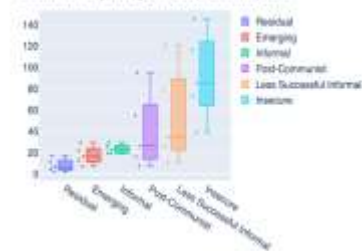
- Real-world decisions rely on multiple charts, not isolated visuals.
- Misinterpreting distributed visuals can lead to wrong conclusions, misleading headlines, and policy errors.
- Headlines often show contradictory or incomplete visuals.
- Small misalignments in axes, time ranges, or metrics can produce misleading narratives.

# Charts are Everywhere, But Can AI Really Understand Them?

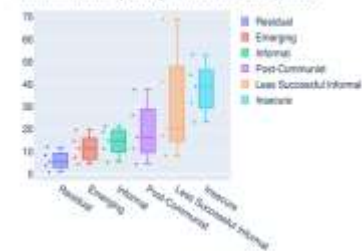
capita (USD) by Country



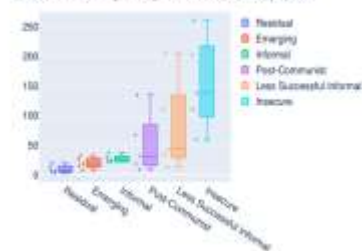
Infant mortality rate (per 1,000 live births), 2004



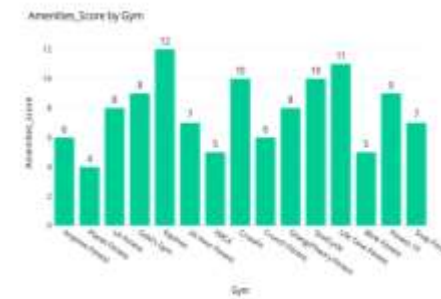
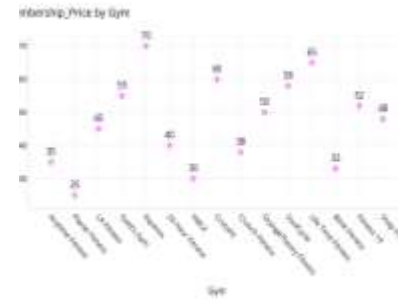
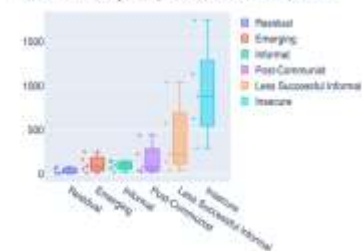
Neonatal mortality rate (per 1,000 live births), 2000



Under-5 mortality rate (per 1,000 live births), 2004



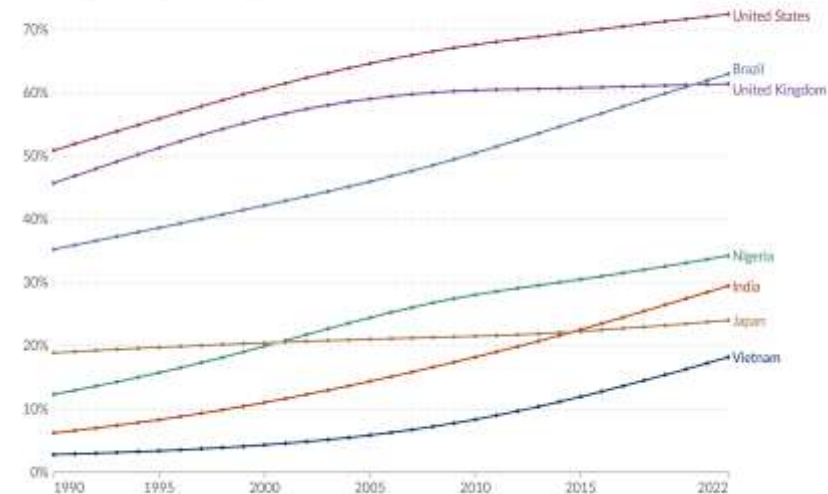
Maternal mortality rate (per 100,000 live births), 2000



## Share of adults who are overweight or obese

"Overweight" is defined here as a body mass index (BMI) above 25. BMI is a person's weight in kilograms divided by their height in meters squared.

Our World in Data



Data source: World Health Organization - Global Health Observatory (2025)

OurWorldinData.org/obesity | CC BY

Note: To allow for comparisons between countries and over time, this metric is age-standardized.

Multi-chart reasoning isn't just vision; it's structured analytical thinking.

# Introducing INTERCHART

**Task:** Evaluate VLMs on distributed, multi-chart understanding.

**Dataset:** Three tiers: DECAF (simple), SPECTRA (relational), STORM (real-world multi-chart).

**Benchmark:** 12k QA pairs, multiple prompting strategies, LLM-as-Judge evaluation.

## Question and Chart Complexity:

- Measures performance across increasing visual + reasoning difficulty:
  - Local lookup → relational comparison → multi-step temporal synthesis
- Tests abilities that current VLMs consistently fail:
  - Trend alignment, multi-chart correlation, cross-entity tracking, time-based reasoning

# Dataset- DECAF

## DECAF - Decomposed Elementary Charts with Answerable Facts

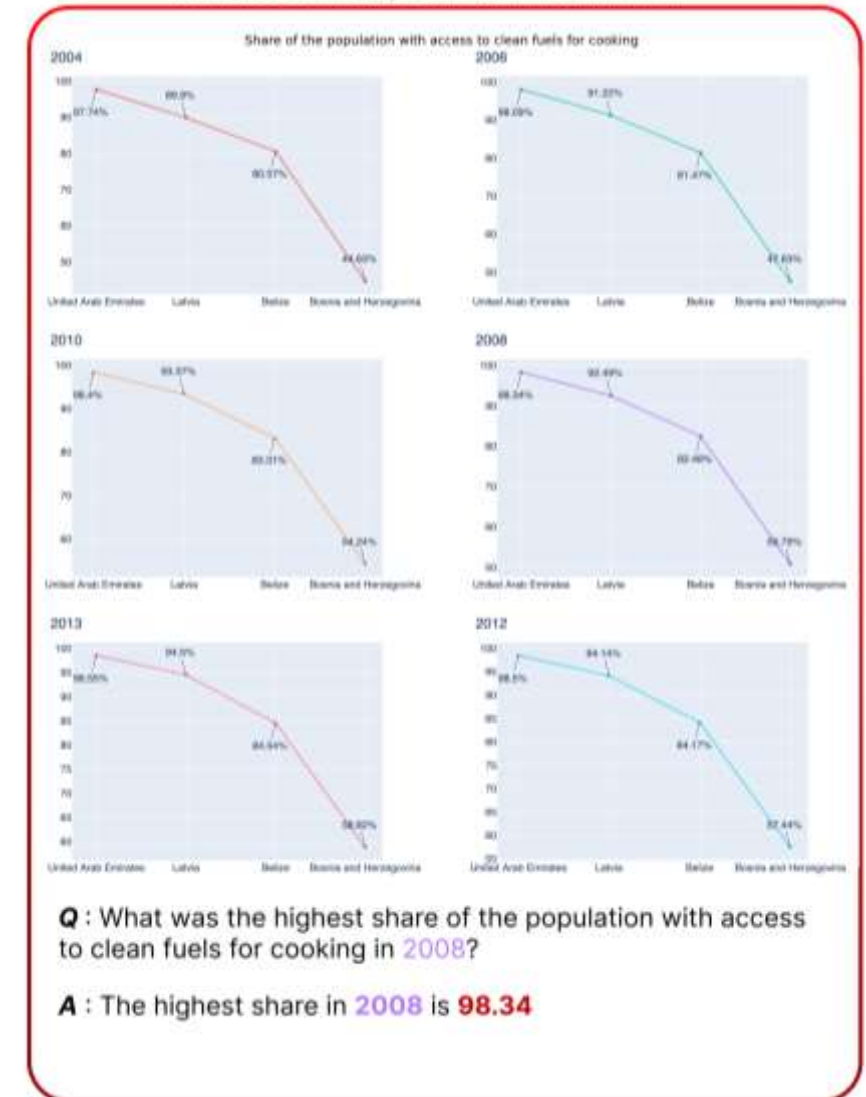
### What DECAF Evaluates

- Basic fact extraction
- Trend identification within a single chart
- Reading axes, labels, and numeric values

### Data Composition

- 355 original charts → 1,188 decomposed charts
- Multiple chart types: line, bar, heatmap, dot, box plot
- 2,809 QA pairs from mixed generation methods (human, LLM, SQL-LLM) Purpose

DECAF provides a controlled baseline for evaluating fundamental visual understanding before introducing relational or temporal complexity.





# Dataset- SPECTRA

SPECTRA - Synthetic Plots for Event-based Correlated Trend Reasoning and Analysis

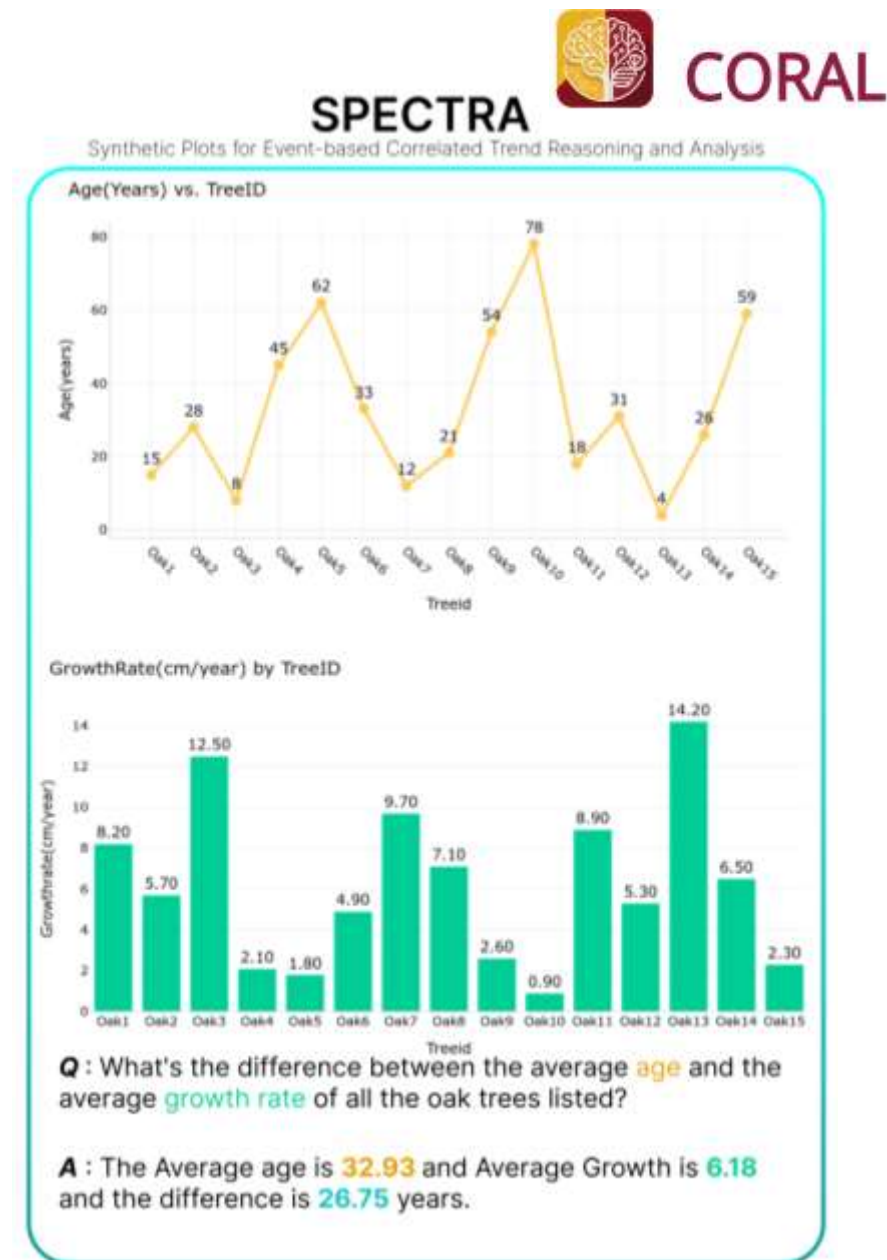
## What SPECTRA Evaluates

- Identifying correlated vs independent trends
- Comparing slopes, averages, and ranges
- Multi-step reasoning across two aligned charts

## Data Composition

- 870 unique charts (paired into 333 context sets)
- 1,717 QA pairs, Includes correlated (1,481) and independent (245) chart pairs

SPECTRA introduces relational complexity, measuring a model's ability to synthesize information across multiple visual sources.



# Dataset- STORM

STORM - Sequential Temporal reasoning Over Real-world Multi-domain charts

## What STORM Evaluates

- Matching events across different timelines
- Inferring numerical ranges and extrema
- Linking economic, demographic, or health indicators across domains
- Multi-step reasoning across visually and semantically distinct charts

## Data Composition

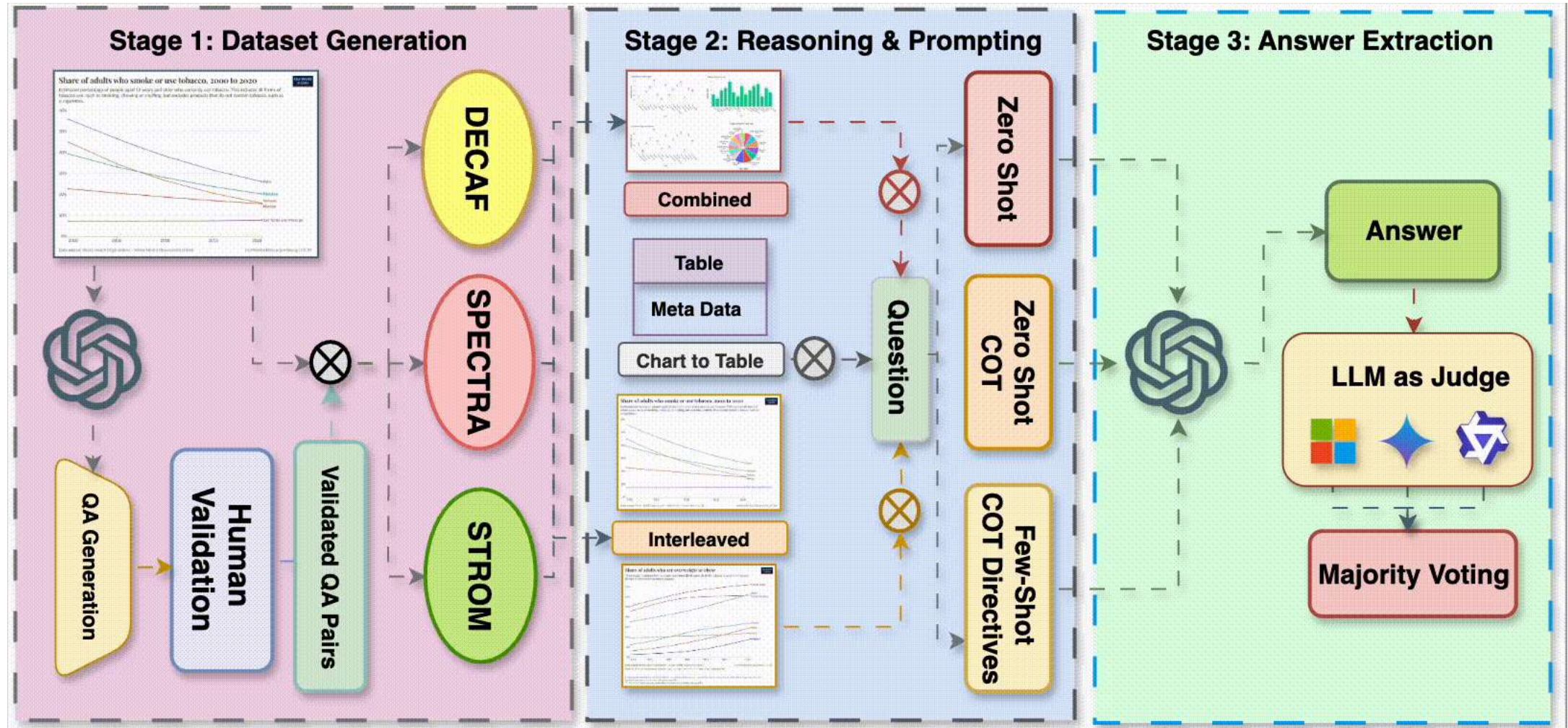
- Data Composition 324 original charts → 648 curated chart images
- 768 QA pairs across three reasoning types:
  - Range Estimation (198)
  - Abstract Numerical Reasoning (275)
  - Entity Inference (295)

STORM represents the highest difficulty tier, mirroring real-world analytical tasks that require integrating distributed, heterogeneous visual information.





# INTERCHART Architecture





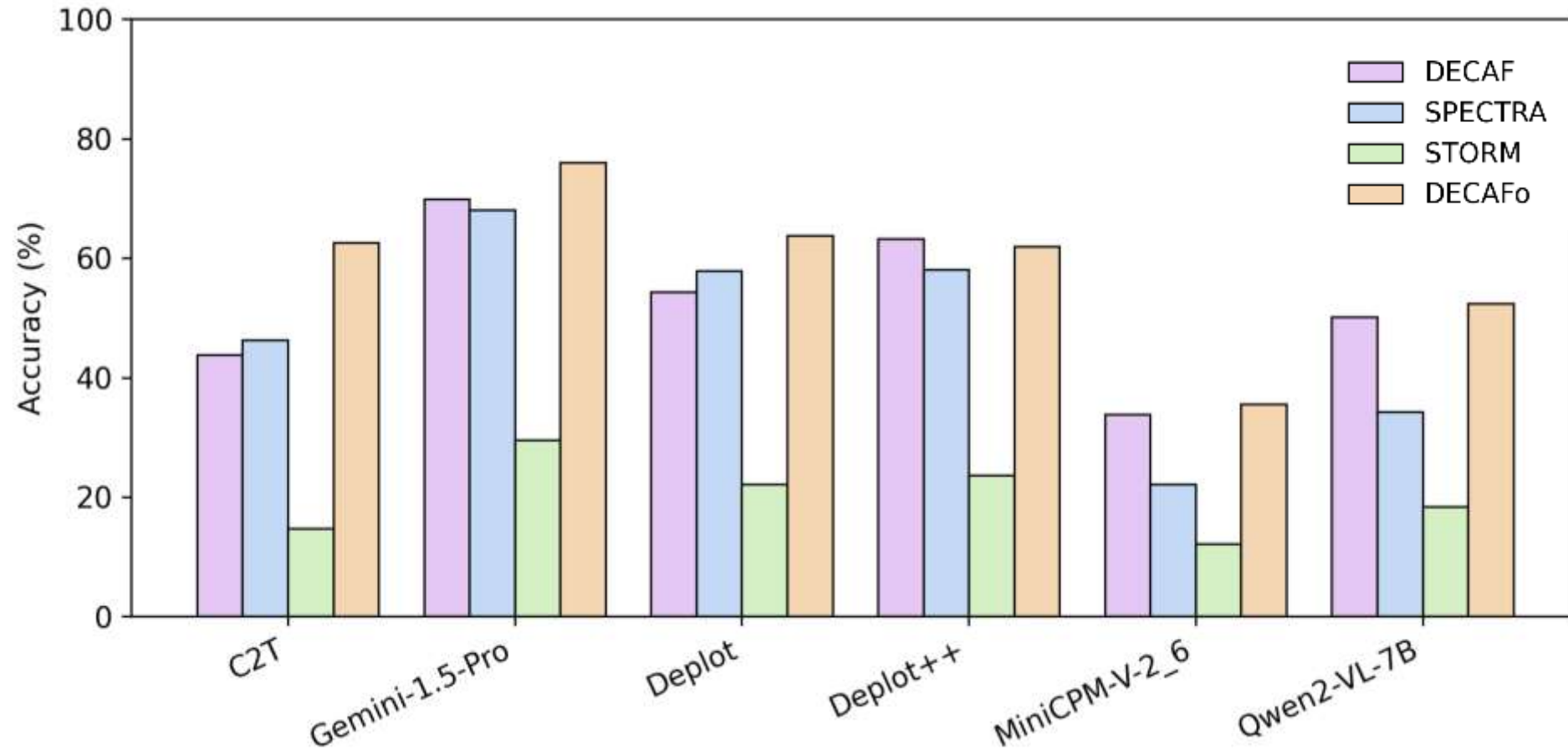
# Results: Interleaved Visual Context



# Results: Combined Visual Context



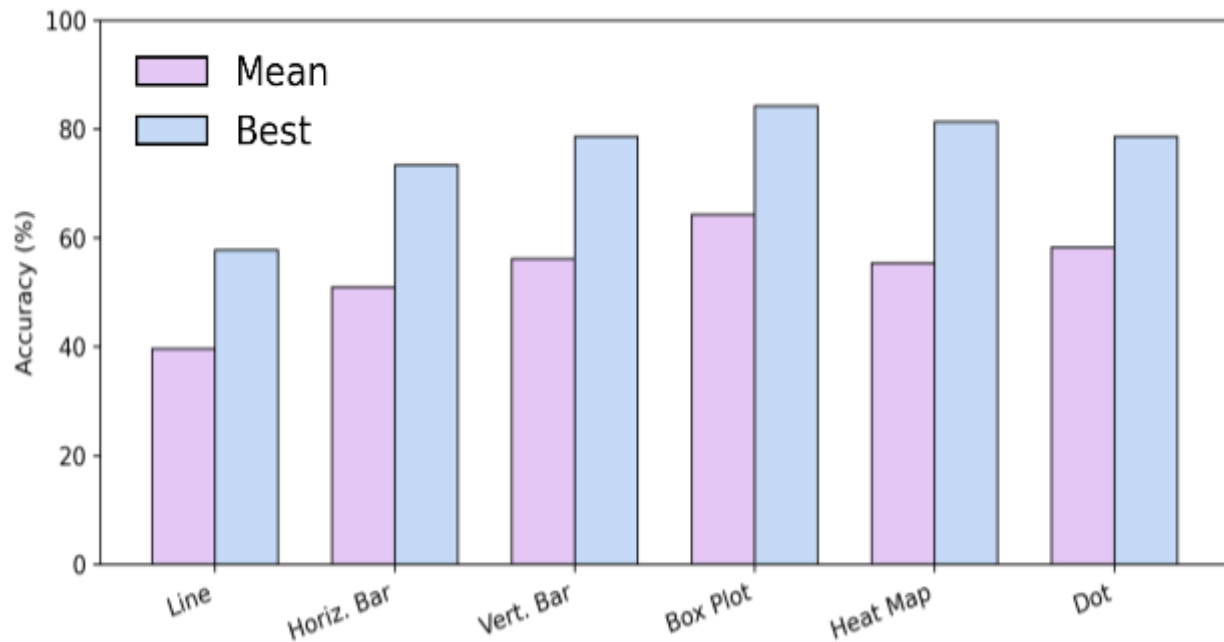
# Results: Chart To Table



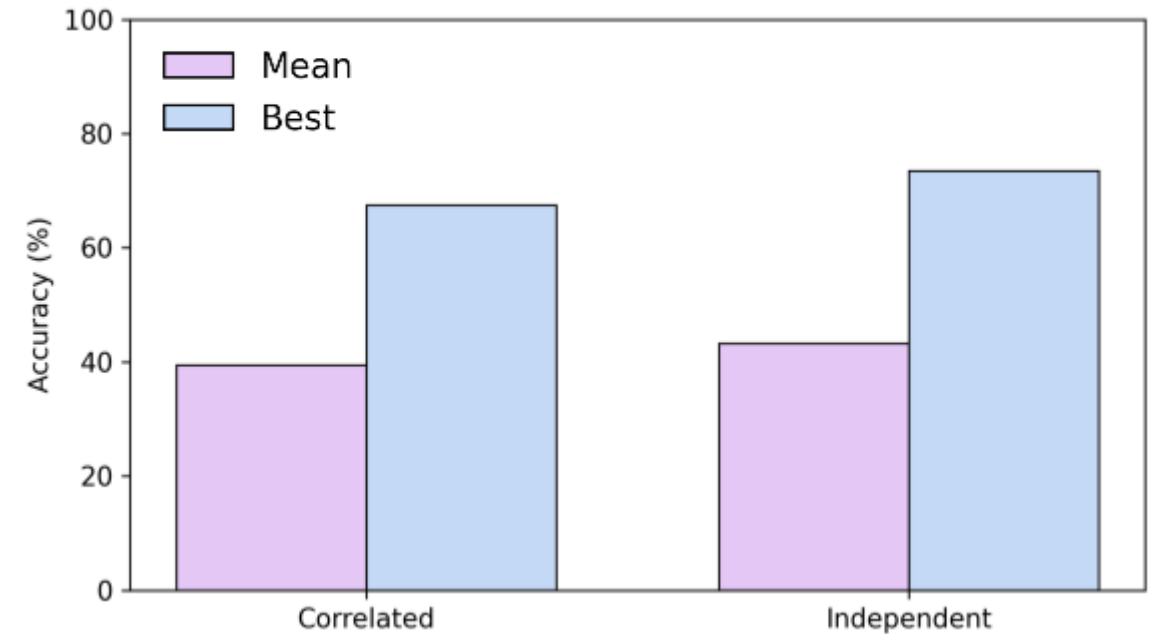


# Results: DECAF and SPECTRA

DECAF by chart type (Mean / Best):

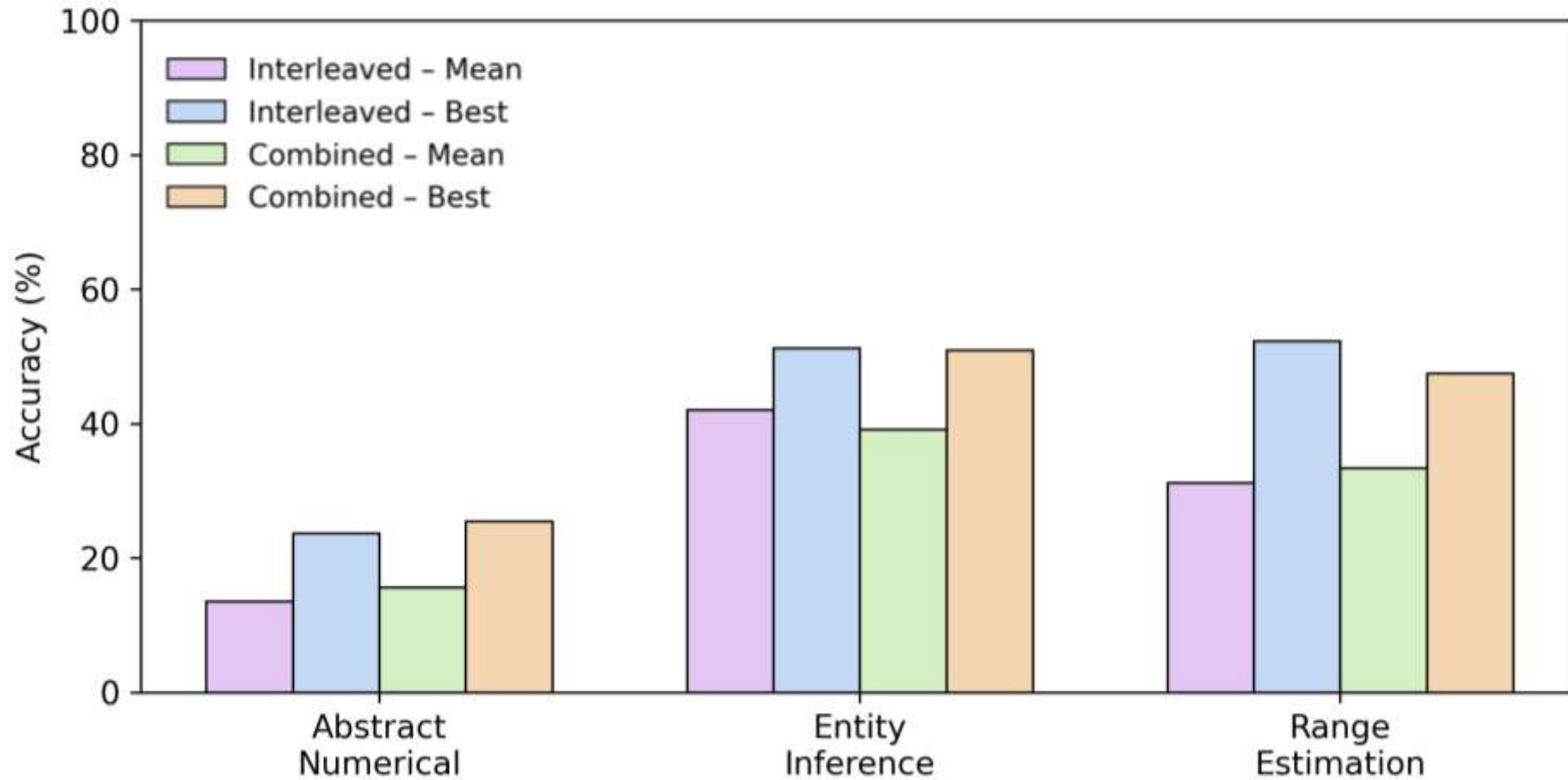


SPECTRA question categories  
(Correlated vs. Independent; Mean / Best).



# Results: STROM

STORM reasoning types (Abstract Numerical, Entity Inference, Range Estimation)  
under Interleaved vs. Combined formats (Mean / Best).



# Comparison in Chart-Based VQA

Dimension	InfoChartQA	ChartMind	ChartQAPro	INTERCHART
<i>Chart Type</i>	Infographics	Mixed	Plots	<b>Plots</b>
<i>Multi-Chart</i>	No	Limited	No	<b>Yes</b>
<i>Real-World Data</i>	Yes	Yes	Yes	<b>Yes</b>
<i>Semantic Drift</i>	Medium	Medium	Low	<b>High</b>
<i>Temporal Reasoning</i>	Low	Medium	Low	<b>High</b>
<i>Visual Diversity</i>	High	High	Low	<b>High</b>
<i>QA Type</i>	Factoid	Hybrid	Factual	<b>Fact + Inference</b>
<i>Evaluation Method</i>	BLEURT	BLEU/LLM	Exact Match	<b>LLM Majority Voting</b>



# Takeaways

1. INTERCHART exposes *systematic failures* in current VLMs, especially when reasoning must integrate information across multiple heterogeneous charts.
2. Decomposing complex visuals into simpler units significantly boosts performance, highlighting that models still rely on localized rather than global reasoning.
3. Accuracy drops sharply from synthetic (SPECTRA) to real-world multi-chart settings (STORM), revealing poor generalization to semantic drift and temporal alignment.
4. Even top-tier VLMs plateau on STORM, emphasizing the need for new architectures that explicitly model cross-chart reasoning, not just visual parsing.

# Thank you

## Team



Anirudh Iyengar



Srija  
Mukhopadhyay



Adnan Qidwai



Shubhankar  
Singh



Dan Roth



Vivek Gupta

Scan for the  
website

